Data exploration Student response addresses the most important characteristics of the dataset and uses these characteristics to inform their decision making. Important characteristics must include: Total number of students: 395 Number of students who passed: 265 Number of students who failed: 130 Number of features: 30 Graduation rate of the class: 67.09% **Preparing the Data** Identify feature and Code has been executed in the iPython notebook, with proper output and no errors. target columns In [4]: # Extract feature (X) and target (y) columns feature\_cols = list(student\_data.columns[:-1]) # all columns but last are features target\_col = student\_data.columns[-1] # last column is the target/label print "Feature column(s):-\n{}".format(feature\_cols) print "Target column: {}".format(target\_col) x\_all = student\_data[feature\_cols] # feature values for all students y\_all = student\_data[target\_col] # corresponding targets/labels print "\nFeature values:-' X\_all.head() # print the first 5 rows Feature column(s):-['school', 'sex', 'age', 'address', 'famsize', 'Pstatus', 'Medu', 'Fedu', 'Mjob', 'Fjob', 'reason', 'guardian', 'trav eltime', 'studytime', 'failures', 'schoolsup', 'famsup', 'paid', 'activities', 'nursery', 'higher', 'internet', 'roma ntic', 'famrel', 'freetime', 'goout', 'Dalc', 'Walc', 'health', 'absences'] Target column: passed Feature values:school sex age address famsize Pstatus Medu Fedu Mjob Fjob \ 0 GP F 18 U GT3 A 4 4 at\_home teacher
1 GP F 17 U GT3 T 1 1 at\_home other
2 GP F 15 U LE3 T 1 1 at\_home other
3 GP F 15 U GT3 T 4 2 health services
4 GP F 16 U GT3 T 3 3 other other ... higher internet romantic famrel freetime goout Dalc Walc health \ ... yes no no 4 3 4 1 1 3 ... yes yes no 5 3 3 1 1 3 yes yes no . . . yes yes yes absences 1 10 3 2 **Preprocess feature** Code has been executed in the iPython notebook, with proper output and no errors. Preprocess teature columns columns As you can see, there are several non-numeric columns that need to be converted! Many of them are simply yes/no, e.g. internet. These can be reasonably converted into 1/0 (binary) values. Other columns, like Mjob and Fjob, have more than two values, and are known as categorical variables. The recommended way to handle such a column is to create as many columns as possible values (e.g. Fjob teacher, Fjob other, Fjob services, etc.), and assign a 1 to one of them and 0 to all others. These generated columns are sometimes called dummy variables, and we will use the pandas.get\_dummies() function to perform this transformation. In [5]: # Preprocess feature columns def preprocess\_features(X): outX = pd.DataFrame(index=X.index) # output dataframe, initially empty # Check each column for col, col data in X.iteritems(): # If data type is non-numeric, try to replace all yes/no values with 1/0 if col\_data.dtype == object: col\_data = col\_data.replace(['yes', 'no'], [1, 0]) # Note: This should change the data type for yes/no columns to int # If still non-numeric, convert to one or more dummy variables if col\_data.dtype == object: col\_data = pd.get\_dummies(col\_data, prefix=col) # e.g. 'school' => 'school\_GP', 'school\_MS' outX = outX.join(col\_data) # collect column(s) in output dataframe return outX X\_all = preprocess\_features(X\_all) print "Processed feature columns ({}):-\n{}".format(len(X\_all.columns), list(X\_all.columns)) Processed feature columns (48):-['school\_GP', 'school\_MS', 'sex\_F', 'sex\_M', 'age', 'address\_R', 'address\_U', 'famsize\_GT3', 'famsize\_LE3', 'Pstatus\_A', 'Pstatus\_T', 'Medu', 'Fedu', 'Mjob\_at\_home', 'Mjob\_health', 'Mjob\_other', 'Mjob\_services', 'Mjob\_teacher', 'Fjob\_ at\_home', 'Fjob\_health', 'Fjob\_other', 'Fjob\_services', 'Fjob\_teacher', 'reason\_course', 'reason\_home', 'reason\_other', 'reason\_reputation', 'guardian\_father', 'guardian\_mother', 'guardian\_other', 'traveltime', 'studytime', 'failures', 'schoolsup', 'famsup', 'paid', 'activities', 'nursery', 'higher', 'internet', 'romantic', 'famrel', 'freetime', 'go out'. 'Dalc'. 'Walc'. 'health'. 'absences'] Split data into Training and test sets have been generated by randomly sampling the overall dataset. training and test Split data into training and test sets sets So far, we have converted all categorical features into numeric values. In this next step, we split the data (both features and corresponding labels) into training and test sets. In [124]: from sklearn.cross\_validation import train\_test\_split # First, decide how many training vs test samples you want num\_all = student\_data.shape[0] # same as len(student\_data) num\_train = 300 # about 75% of the data num\_test = num\_all - num\_train # TODO: Then, select features (X) and corresponding labels (y) for the training and test sets # Note: Shuffle the data or randomly select samples to avoid any bias due to ordering in the dataset # I am using treain\_test\_split which wraps input validation and next(iter(ShuffleSplit(n\_samples))) and application to X\_train , X\_test, y\_train , y\_test = train\_test\_split(X\_all, y\_all, test\_size = num\_test, random\_state=0 ) print "Training set: {} samples".format(X train.shape[0]) print "Test set: {} samples".format(X\_test.shape[0]) # Note: If you need a validation set, extract it from within training data Training set: 300 samples Test set: 95 samples **Training and Evaluating Models** Time and Space Both the big-O notation for the space complexity to represent the model and the time for the Complexity algorithm to make a prediction are provided, or a list of several of the major factors that affect the time & space complexities are presented with a description of the largest driving factor as constant, linear, logarithmic, polynomial, etc in nature. Student presents resources or reasonable justification for their response. Due to small dataset the prediction and training times of the model are negleble. Driving factor of comuputionn is liner in nature, meaning is we increase number of samples, the time it will take to execute will increase. As such driving factor is Liner.https://en.wikipedia.org/wiki/Time\_complexity#Polynomial\_time **Model Application** The pros and cons or application for each model is provided with reasonable justification why each model was chosen to explore. **Decision Tree** Classifier • What are the general applications of this model? Decision Tree Classifier applies a straitforward idea to solve the classification problem. it poses a series of carefully crafted questions about the attributes of the test record. Each time time it receive an answer, a follow-up question is asked until a conclusion about the calss label of the record is reached. What are its strengths and weaknesses? Advantages: Decision trees are simple to use, easy to understand. Disadvantages: Even a small change in input data can at times, cause large changes in the tree. Decision trees are also prone to errors in classification, owing to differences in perceptions and the limitations of applying statistical tools. Given what you know about the data so far, why did you choose this model to apply? Predictive Power, Relative Simplicity The pros and cons or application for each model is provided with reasonable justification why each Model Application model was chosen to explore. **KNeighborsClassifier** • What are the general applications of this model? KNeighborsClassifier classifies data by evaluating classes and data point and classes of K nearest neighbors. Class is determined based on majority vote. In case of student data we only have classes and we measure distance to nearest neighbor (trough parameter search process we found that number to be 7). When the model encounters an example it has not see it will compare the distance of such example from neighbors around it. Model uses neighbors majority vote to determine class of the new example. Its important to use uneven number of neighbors to achieve tie breaks in majority voting. What are its strengths and weaknesses? Advantages: Relatively fast as only distances are measured. • Robust to noisy training data Disadvantages: Initial placement of cluster centroid can affect clustering, we have to adjust centroids as we cluster. Even after this examples that are close together can switch cluster assignments. • Need to determine value of parameter K (number of nearest neighbors) • Given what you know about the data so far, why did you choose this model to apply? I originally chose this model as one of the possible models because it does well generalizing data and selecting membership based on similarity. Meaning if new evaluation point is close or falls into an existing group then we can easily clarify it based on distance from other members of its class. The pros and cons or application for each model is provided with reasonable justification why each Model Application model was chosen to explore. SVM What are the general applications of this model? SVM classifies data by separating classes and than placing as wide of a margin as possible that will still achieve linear separation. We can use this model to linearly separate the data with little fear that outliers will cause imbalance. What are its strengths and weaknesses? • **Advantage** SVMs have good generalization performance

**Disadvantage:** limitation of the support vector approach lies in choice of the kernel Burgess (1998)

All the required time and F1 scores for each model and training set sizes are provided within the

Training time (secs) Prediction time (secs)

F1 score for training set

F1 score for test set

Training time (secs)

Prediction time (secs)

F1 score for test set

Training time (secs)

Prediction time (secs)

F1 score for test set

• Based on the experiments you performed earlier, in 1-2 paragraphs explain to the board of supervisors

what single model you chose as the best model. Which model is generally the most appropriate based on

Based on tests performed, I compared 3 models, Decision tree clasifier, Random forest, and Support

Vector Classification. After evaluating performance of each model I choose Decision Tree Classifier

F1 score for training set | 0.781

F1 score for training set | 0.561 | 0.540

allowing for classes of new data points to be determined with limited overfitting.

chart given. The performance metrics are reasonable relative to other models measured.

I originally chose this model as one of the possible models because it separates data with a margin

Training set size

200

0.830

0.802 | 0.794 | 0.785

Training set size

200

0.007 | 0.018 | 0.024

0.444 | 0.547 | 0.768

Training set size

200

0.001

0.002

0.830

0.783

100

0.001

0.002

0.768

300

0.829

300

0.833

0.001 0.006

0.841

0.788

0

100

0.805

100

Given what you know about the data so far, why did you choose this model to apply?

Student Intervention System - Decision Tree Classifier

Student Intervention System - LinearSVC Classifier

Student Intervention System - KNeighborsClassifier Classifier

the available data, limited resources, cost, and performance?

Model for number of reasons:

and the recall.

a prediction).

**Classification vs Regression** 

based on avialble data

Which type of supervised machine learning problem is this, classification or regression? Why?

This is a classification problem becasuse we are trying to predict if student will drop out or not

Classification vs

**Exploring the Data** 

Regression

## Model Performance **Metrics**

**Choosing the Best Model** 

Choosing the

Optimal Model

**Terms** 

example Describing the Model in Layman's **Model Tuning** The final model chosen is fined tuned using at least one parameter with at least three settings. I have tuned each model to find best parmaeter. Once the models have been tuned following are f1 scores on test data: F1 Messure: 0.76 [DecisionTreeClassifier] F1 Messure: 0.68 [KNeighborsClassifier]

What is the model's final F<sub>1</sub> score? Best F1 score for test set is 0.788 Before model tuning following were the F1 scores on the training data. F1 Messure: 0.71 [DecisionTreeClassifier] F1 Messure: 0.77 [KNeighborsClassifier] F1 Messure: 0.81 [SVC] After mode tuning F1 Messure: 0.80 [DecisionTreeClassifier] F1 Messure: 0.80 [KNeighborsClassifier] F1 Messure: 0.80 [SVC]

a

F1 Messure: 0.71 [SVC]

Best F1 score for test set is 0.788

Code reflects the description in the documentation.

model built.

**Tuned F1 Score** 

**Quality of Code** 

**Functionality** 

KNeighborsClassifier Classifier uses a majority neibour voting to detirmine the class of a new Fine-tune the model. Use Gridsearch with at least one important parameter tuned and with at least 3 settings. Use the entire training set for this. Without finding the best parameters following were the prediction accurancies: I have tuned each model to find best parmaeter. Please refer to outputs of each model KNeighborsClassifier Classifier uses a majority neibour voting to detirmine the class of a new example. For example if we have 2 classes 'a' and 'o' and we need to find a class for new value of 'c'. We will examine chosen number of neighbors (k = 3) to determine the class of 'c' by taking a majority vote. In this case we have 2 'o' and 1 'a' so the new class is detrained to be majority 'o'. Given N training vectors, kNN algorithm identifies the k nearest neighbors of 'c', regardless of labels

I have used grid search to search for the best parameters to use in the model. Grid

The F1 score is provided from the tuned model and performs better than the default model chosen.

search was done on all available data as only the parameters are being chosen, vs

 Accurancy: Model predicts with high F1 Score. F1 score conveys the balance between the precision - Peformance: Obtainings training and test prediction fastes

In 1-2 paragraphs explain to the board of supervisors in layman's terms how the final model chosen is supposed to work (for example if you chose a Decision Tree or Support Vector Machine, how does it make Example k = 3
 classes 'a' and 'o'
 find class for 'c'